# Overview of cluster management tools

Open Science Grid

THE UNIVERSITY OF
CHICAGO

Marco Mambelli – marco@hep.uchicago.edu

August 9 2011

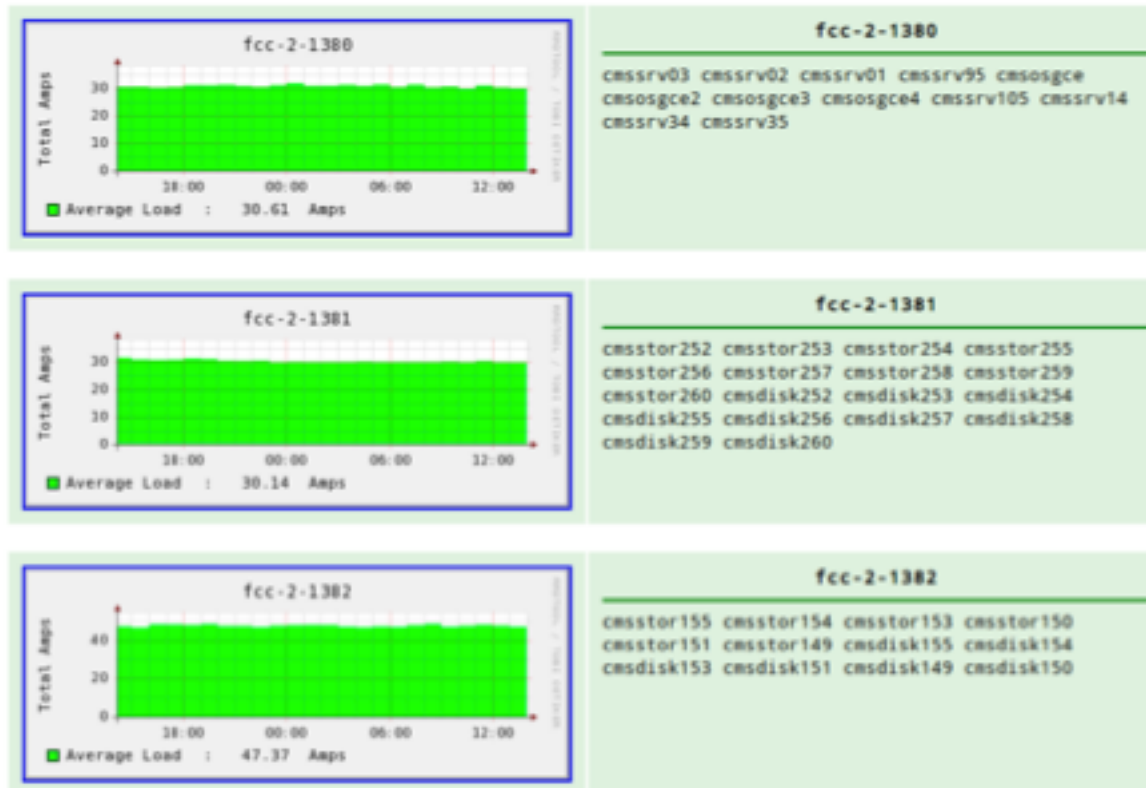OSG Summer Workshop

TTU - Lubbock, TX

# Cluster Management Overview

- Management infrastructure
- Provisioning
- Configuration and software package management
- Monitoring

# Management Infrastructure

- Remote power-cycling and serial console access

- FEF has standardized on Avocent ACS serial console servers and Avocent PM line of PDU

- Other depts use IPMI or a mix of Avocent and APC products

- All depts have scripts to control and configure remote power- cycling and serial console access

# CMS Tier 1 Power Usage Plots



**fcc-2-1380**

cmssrv03 cmssrv02 cmssrv01 cmssrv95 cmsosgce
cmsosgce2 cmsosgce3 cmsosgce4 cmssrv105 cmssrv14
cmssrv34 cmssrv35

**fcc-2-1381**

cmsstor252 cmsstor253 cmsstor254 cmsstor255
cmsstor256 cmsstor257 cmsstor258 cmsstor259
cmsstor260 cmsdisk252 cmsdisk253 cmsdisk254
cmsdisk255 cmsdisk256 cmsdisk257 cmsdisk258
cmsdisk259 cmsdisk260

**fcc-2-1382**

cmsstor155 cmsstor154 cmsstor153 cmsstor150
cmsstor151 cmsstor149 cmsdisk155 cmsdisk154
cmsdisk153 cmsdisk151 cmsdisk149 cmsdisk150

CMS plots power utilization by querying PDUs using SNMP. This data can be particularly useful to datacenter managers

Open Science Grid

# Provisioning Tools

- ## Provisioning Tools

  - Preparing a system for use; OS installation and initial configuration

- ## Tools

  - PXE/Kickstart

  - Rocks

  - Cobbler

  - Perceus

# FEF's PXE/Kickstart Setup

- FEF uses custom tool built on top of MySQL and Perl DCHP server modules

- No dhcpd restart required

- Web front-end for specifying kickstart / nodecombinations

- Very flexible

- Kickstart files are created dynamically based on selections from the web GUI

- Planning to eval Cobbler later this year.

# Rocks Clusters

- Open source Linux distribution based on CentOS

- Created in 2000

- Created for easy deployment of large clusters

- Used by CMS Tier 1 at Fermilab for 2400 machines

- CMS can install ~500 systems in 1 hour

- Only one OS version per Rocks server may be a deal breaker for some

# Cobbler

- Cobbler is infrastructure to provision a node
- Cobbler is RedHat specific, uses kickstart
- Cobbler holds OS, profiles and system data
- A new system requires a profile, a MAC address and a name, nothing more.
- Provisions new system via PXE boot, via settings for individual system.
- CLI and Web GUI

# Cobler at MWT2

- Cobbler is implemented through multiple services on a single server. It acts as a TFTP server for PXE booting, controls the repositories used for install, and provides DHCP/DNS support as well.

# Configuration and Package Management

- Tools that help manage system configuration (files, dirs, permissions, etc.) and software packages

- Popular open source solutions
  - Cfengine
  - Puppet
  - Bcfg2
  - Warewulf

# Cfengine

- **First version released in 1993**
- **Written in C**
- **Fairly easy to understand syntax**
- **Relatively easy to find sysadmins with experience**
- **FEF and UWM used for many years**

# Puppet

- New generation of configuration management system

- Extensible, declarative language

- Understands dependencies (huge benefit)

- Better reporting than Cfengine

- Auto generation of documentation (think Javadoc)

# FEF Puppet Usage

- Management of all external mounts
- Kerberos files -- keytab files, .k5login, etc
- Package management (RPM sets grouped by cluster)
- NIC bonding configs
- Group quotas
- Grid host certs
- FEF_backup
- FEF avg is 325 actions per node every Puppet run

# Cfengine vs. Puppet (High Level)

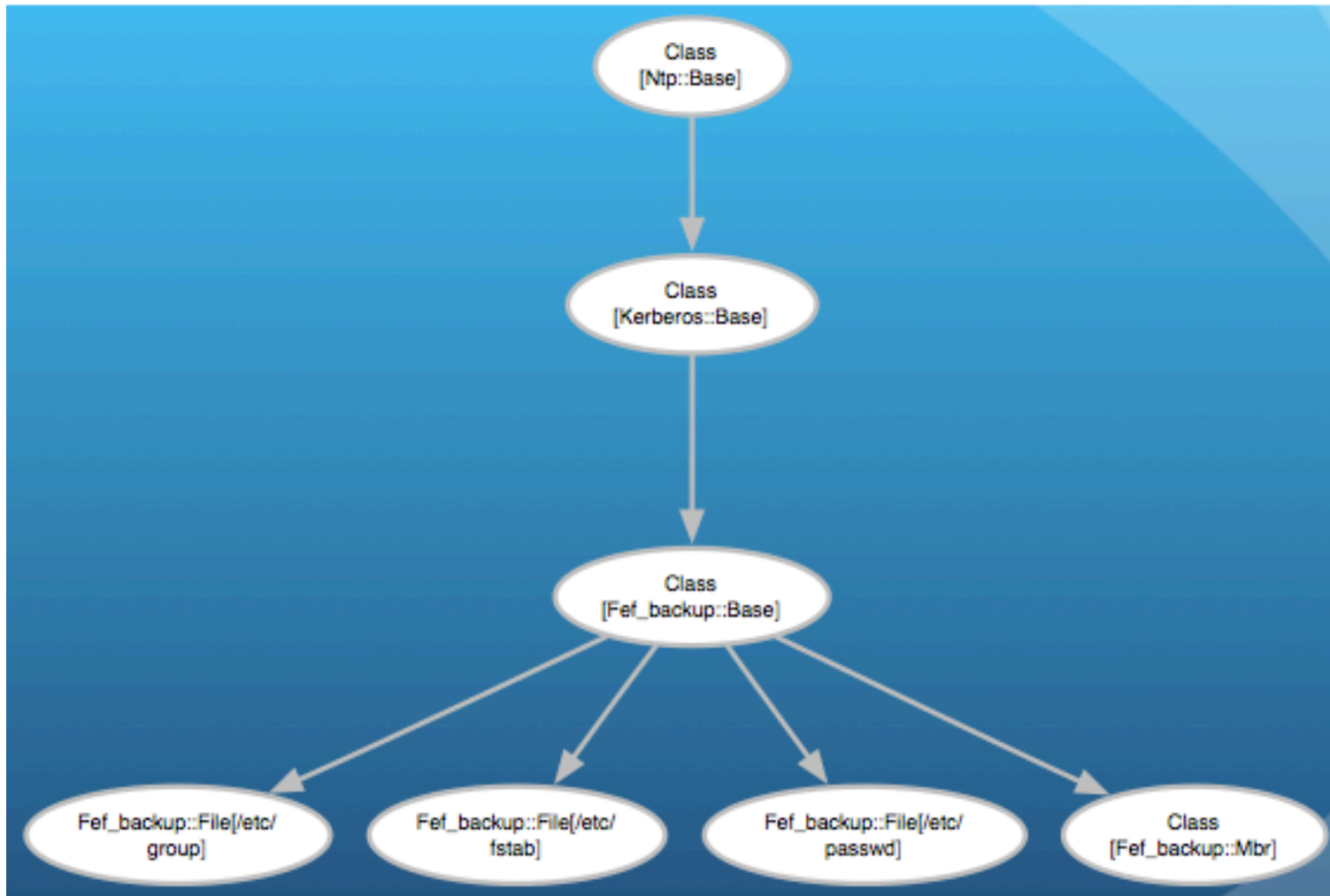| | Native File Editing | Dependency Management | Commercial Support | Dependency Graphs | Scalability |
|---|---|---|---|---|---|
| **Cfengine** | ✔ | | ✔ | | ✔ |
| **Puppet** | | ✔ | ✔ | ✔ | ✔ |

# Puppet Add User Example

```
User { managehome    => true,
       ensure        => present,
       gid           => users,
       shell         => "/bin/bash",
}

user { "mark":
    uid => 1000,
}

user { "fred":
    uid  => 1001,
}

user { "jane":
    uid => 1002,
}
```
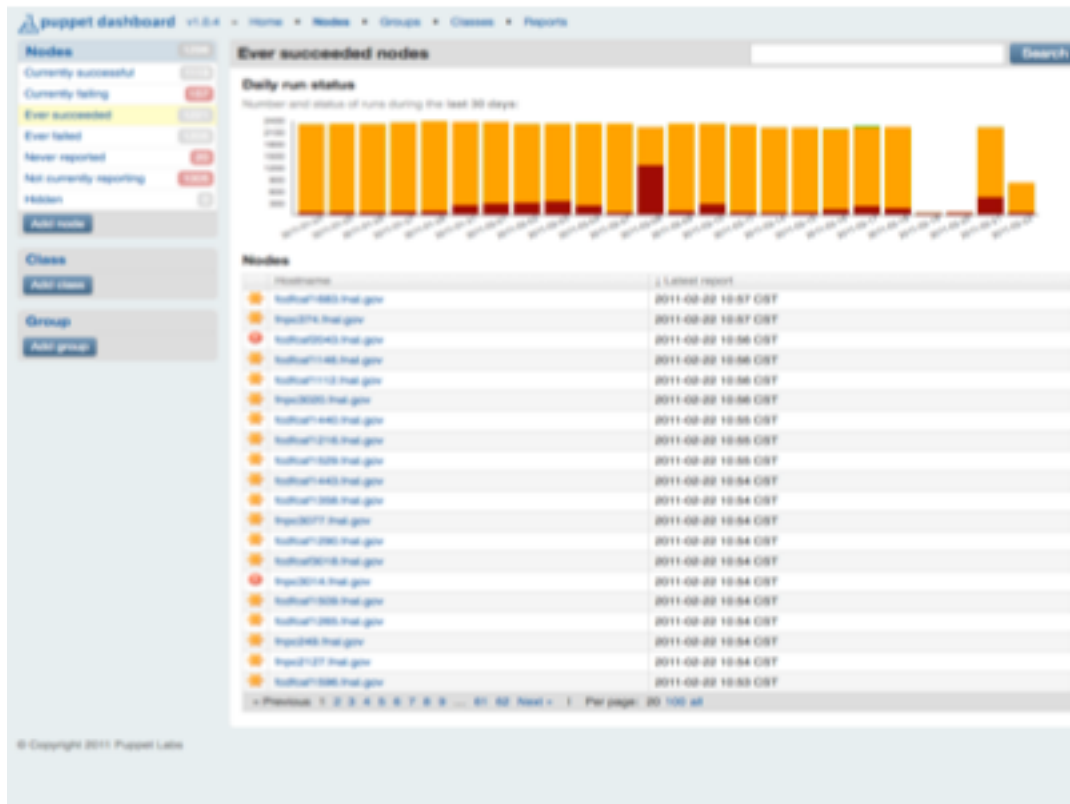
Open Science Grid

# Cfengine Add User Example

```
"pw[mark]" string => "mark:x:1000:100:Mark Burgess:/home/
mark:/bin/bash";

"pw[fred]" string => "fred:x:1001:100:Right Said:/home/
fred:/bin/bash";

"pw[jane]" string => "jane:x:1002:100:Jane Doe:/home/
jane:/bin/bash";

"users" slist => getindices("pw");

files:

  "/etc/passwd"

      edit_line => append_users_starting("addusers.pw");

  "/etc/group"

        edit_line => append_user_field("root","4","@
(addusers.users)");

  "/home/$(users)/."

    create => "true",

      perms => mog("755","$(users)","users");
```

# Example Puppet Dependency Graph

Open Science Grid

# Puppet Dashboard



Puppet dashboard is a web interface for quickly viewing puppet run status and the state of individual system configurations.

Open Science Grid

# Bcfg2

- BCFG2 is an xml-based configuration management system
- Developed by Argonne National Lab
- Being used by CMS Tier 1 at Fermilab for 2 years to manage a limited number of configuration items. Not RPMS
- Developers are very responsive; provide support via mailing list and IRC
- Complex file manipulation can be tricky
- Does simple pre/post dependencies

# Bcfg2 Reporting System



Node status and last run times are viewable from the Bcfg2 web interface

# Monitoring

- **Tools to help monitor system state and performance**
- **In use at Fermilab:**
  - Zabbix
  - Nagios
  - Ganglia
  - Many custom solutions using MRTG, RRDtool, etc

# Zabbix

- Being used by CMS Tier 1 to monitor approx 2.4K nodes; performs 100K checks.

- Does status and performance monitoring.

- Relatively new compared to Nagios.

- Most configuration is done via the web interface.

- Easy to add custom checks and alerts.

Open Science Grid

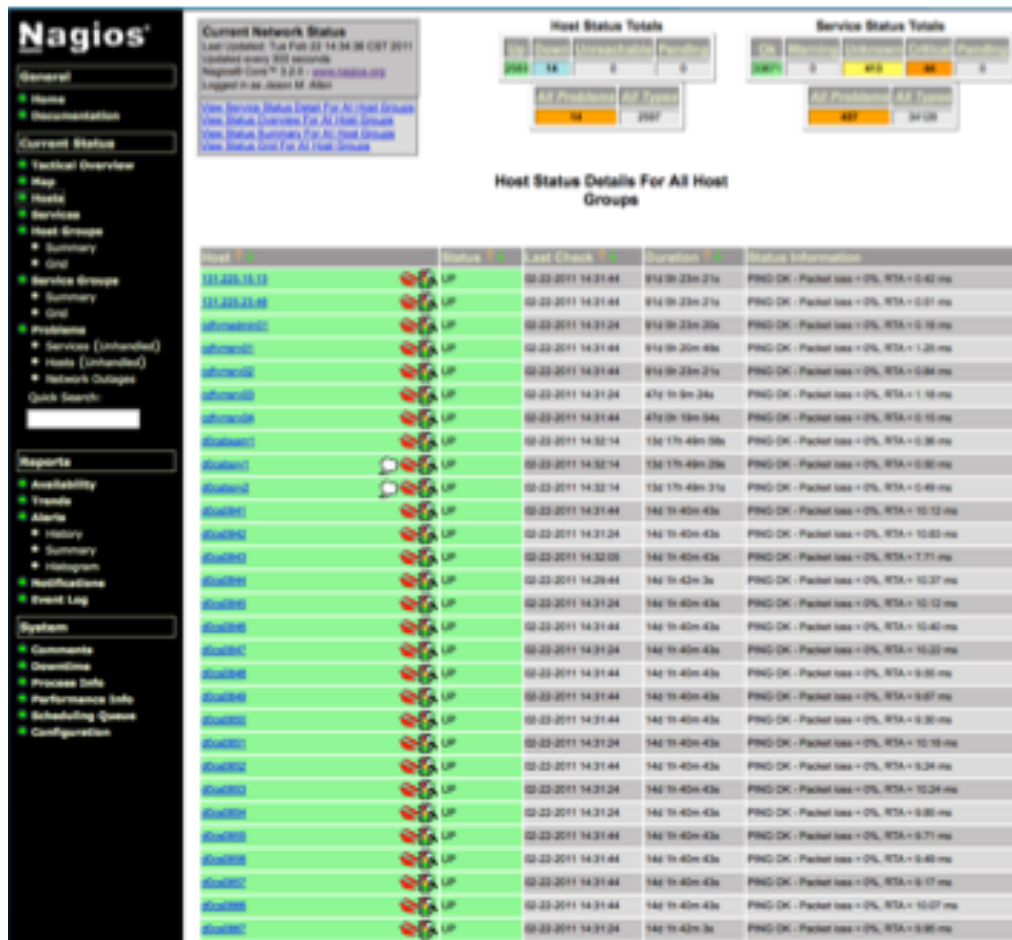# Zabbix Dashboard



Zabbix provides a polished web interface that displays finely grained status and performance information

Open Science Grid

# Nagios

- Used by FEF; 3.5K nodes, approx 30K checks on one server

- Around for many years.

- Create a new check by dropping shell script on node

- (check_mk plugin)

- Nagios support built-in to Puppet.

- Web interface can be slow and feels dated.

# Nagios Web Interface



The Nagios web interface is functional but feels dated. Performance is an issue when monitoring many hosts

Open Science Grid

# OSG Survey

- The most used cluster management tool is Rocks, sometime customized

- Followed by Puppets, Cobbler and Cfengine.

- Users are happy with the performance of the tools they use, especially Rocks and Puppets.

- The difficulty of the first installation is average (sometime long or with some guesswork) to easy (works out of the box); same for the updates.

# OSG Survey (cont)

- The operation is automatic or requires simple documented tasks.

- Rocks seem the easiest to operate

- Puppet is the easiest to install/update

- Available documentation is good for Rocks, good to average (there could be more or sometime is confusing) for the others.

- There are some long time users of Rocks and Cfengine while Puppet gained popularity in recent times.

- Comparison of Puppet/Cfengine/Bcfg2
  - https://cd-docdb.fnal.gov:440/cgi-bin/ShowDocument?docid=3967
- Evaluation and feature comparison of the Nagios and Zabbix monitoring systems
  - http://cd-docdb.fnal.gov/cgi-bin/ShowDocument?docid=3277
- Survey about the setup of clusters in OSG
  - Available from OSG DOCDB

# Credits

- Thank you to Jason Allen, Head of Fermilab Experiments Facilities (FEF) Department in the Computing Division – this talk is based heavily on his Cluster Management talk at the 2011 OSG all hands meeting

?

!

Overview of cluster management tools